

Multi-Target and Structured-Output Learning

INF581 Advanced Machine Learning and Autonomous Agents

Jesse Read



Last Updated: January 12, 2022

Multi-Label Classification

- 1 Multi-Label Classification
- 2 Classifier Chains
- 3 Inference as Search
- 4 Benefits of Structure
- 5 Structure from Data / Learning Structure as a Search
- 6 Summary

Multi-Label Classification

Input	Beach	Sunset	Foliage	Urban
	1	0	1	0
	0	1	0	0
	0	1	0	1
	0	1	1	0
	0	0	1	1
	?	?	?	?

Given an instance \mathbf{x} , we obtain predictions $\hat{\mathbf{y}} = h(\mathbf{x})$.



The Lord of the Rings: The Fellowship of the Ring (2001)



Top 500

PG-13 | 178 min.

Adventure, Fantasy

19 December 2001 (USA)



Your rating: ★ ★ ★ ★ ★ ★ ★ ★ ★ ★ /10

Ratings: 8.8/10 from 1,110,948 users Metascore: 92/100

Reviews: 4,988 user | 294 critic | 34 from Metacritic.com

A meek hobbit of the Shire and eight companions set out on a journey to Mount Doom to destroy the One Ring and the dark lord Sauron.

Director: Peter Jackson

Writers: J.R.R. Tolkien (novel), Fran Walsh (screenplay), [2 more credits »](#)

Stars: Elijah Wood, Ian McKellen, Orlando Bloom |

[See full cast and crew »](#)

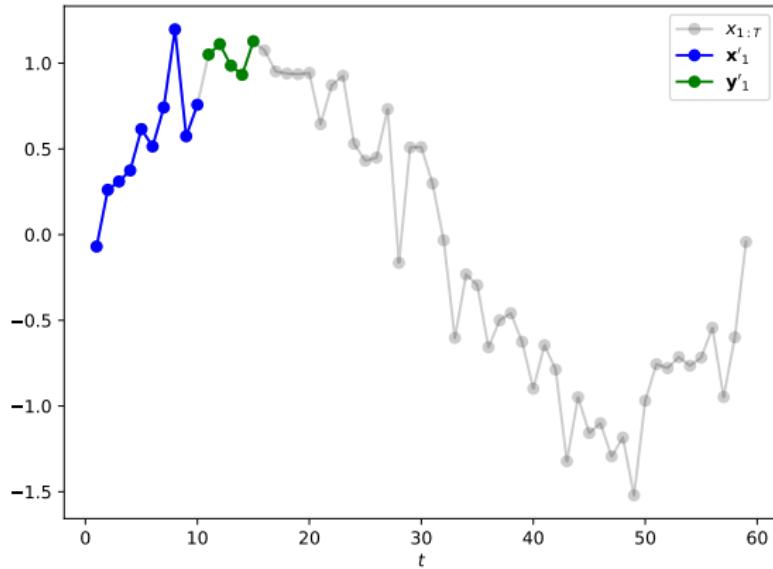


Source: [\[1\]](#)

Missing-data Imputation / Recommender Systems

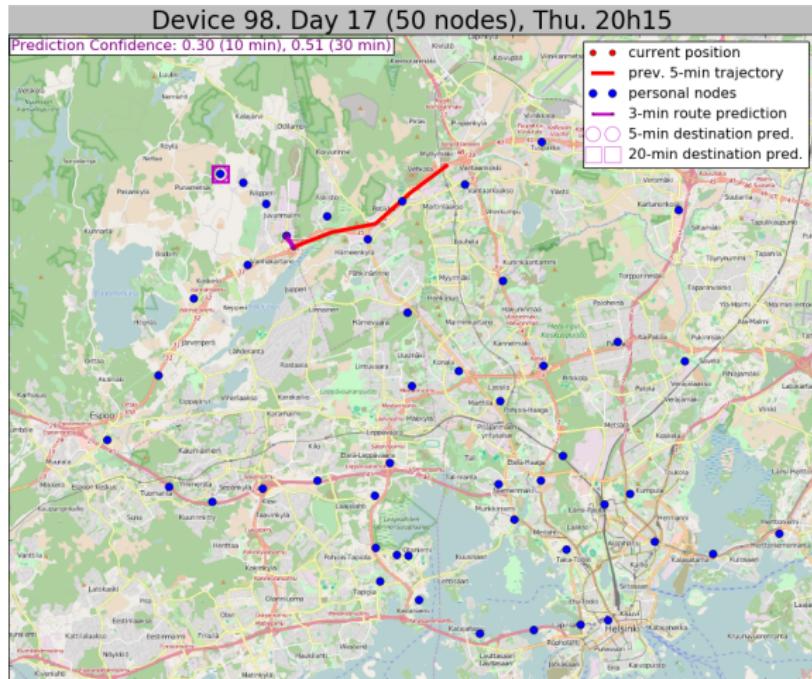
	Film 2 X_2	Film 3 X_4	Book 1 X_1	Book 2 X_3	Song 5 X_5
>User 1 (Red)	0	0	1	1	0
User 2 (Green)	1	1	?	0	?
User 3 (Yellow)	0	0	1	0	0
User 4 (Orange)	1	1	?	0	1
User 5 (Black)	0	0	0	?	?
User 6 (Pink)	1	0	?	1	?

Time Series Forecasting



(including multi-dimensional time series; generalization to $\mathbf{y} \in \mathbb{R}^m$).

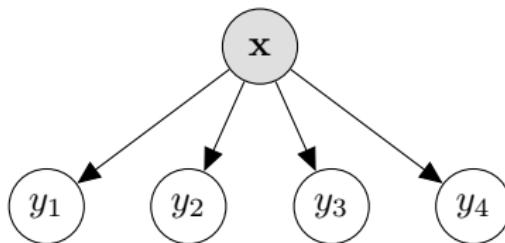
Trajectory Prediction



Trajectory prediction in urban environment using mobile phone data

Binary Relevance (Independent Classifiers): The Baseline

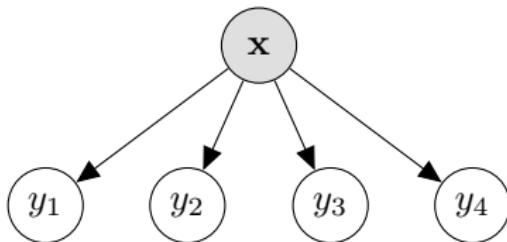
\mathbf{X}	Y_1	Y_2	Y_3	Y_4
$x^{(1)}$	0	1	1	0
$x^{(2)}$	1	0	0	0
$x^{(3)}$	0	1	0	0
$x^{(4)}$	1	0	0	1
$x^{(5)}$	0	0	0	1
$\tilde{\mathbf{x}}$?	?	?	?



The **binary relevance method** = one binary classifier trained for each label, i.e., **independent models**.

Binary Relevance (Independent Classifiers): The Baseline

\mathbf{X}	Y_1	\mathbf{X}	Y_2	\mathbf{X}	Y_3	\mathbf{X}	Y_4
$x^{(1)}$	0	$x^{(1)}$	1	$x^{(1)}$	0	$x^{(1)}$	1
$x^{(2)}$	1	$x^{(2)}$	0	$x^{(2)}$	1	$x^{(2)}$	0
$x^{(3)}$	0	$x^{(3)}$	1	$x^{(3)}$	0	$x^{(3)}$	1
$x^{(4)}$	1	$x^{(4)}$	0	$x^{(4)}$	1	$x^{(4)}$	0
$x^{(5)}$	0	$x^{(5)}$	0	$x^{(5)}$	0	$x^{(5)}$	0
\tilde{x}	?	\tilde{x}	?	\tilde{x}	?	\tilde{x}	?



The **binary relevance method** = one binary classifier trained for each label, i.e., **independent models**.

Main Questions (for this Lecture)

When is **binary relevance** (independent models) *not* the best approach to take?

And in that case, when independent models are not optimal,

- Which **structure**, how to **learn it**
- How to performing **inference** in that structure

And why does this work better than independent models?

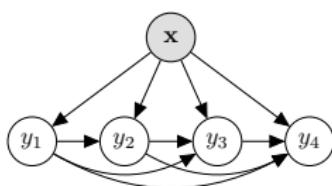
'Bonus question': Overlap between probabilistic and deep neural models

Classifier Chains

- 1 Multi-Label Classification
- 2 Classifier Chains
- 3 Inference as Search
- 4 Benefits of Structure
- 5 Structure from Data / Learning Structure as a Search
- 6 Summary

Classifier Chains (and Greedy Inference)

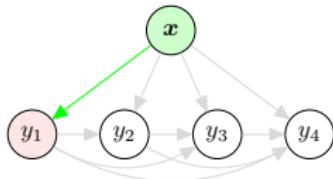
A **classifier chain** is a Bayesian network, with a fully-cascaded structure across output variables. Each $P(Y_j|\text{pa}(Y_j))$ is provided by a **base classifier** (e.g., logistic regression, decision trees, neural network, . . .). Motivation: Model **label dependence** with structure.



X	\hat{y}_1	\hat{y}_2	\hat{y}_3	\hat{y}_4
$x^{(1)}$	0	1	1	1
$x^{(2)}$	1	0	0	0
$x^{(3)}$	0	1	0	1
$x^{(4)}$	1	0	0	0
$x^{(5)}$	0	0	0	0

Classifier Chains (and Greedy Inference)

A **classifier chain** is a Bayesian network, with a fully-cascaded structure across output variables. Each $P(Y_j|\text{pa}(Y_j))$ is provided by a **base classifier** (e.g., logistic regression, decision trees, neural network, . . .). Motivation: Model **label dependence** with structure.

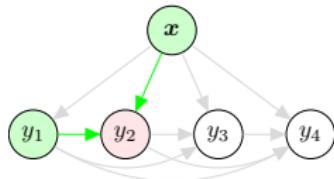


X	Y ₁	Y ₂	Y ₃	Y ₄
$x^{(1)}$	0	1	1	1
$x^{(2)}$	1	0	0	0
$x^{(3)}$	0	1	0	1
$x^{(4)}$	1	0	0	0
$x^{(5)}$	0	0	0	0

\tilde{x}	\hat{y}_1
-------------	-------------

Classifier Chains (and Greedy Inference)

A **classifier chain** is a Bayesian network, with a fully-cascaded structure across output variables. Each $P(Y_j|\text{pa}(Y_j))$ is provided by a **base classifier** (e.g., logistic regression, decision trees, neural network, . . .). Motivation: Model **label dependence** with structure.

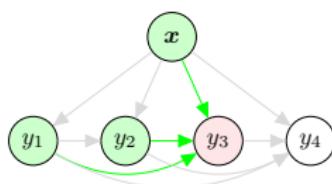


X	Y ₁	Y ₂	Y ₃	Y ₄
$x^{(1)}$	0	1	1	1
$x^{(2)}$	1	0	0	0
$x^{(3)}$	0	1	0	1
$x^{(4)}$	1	0	0	0
$x^{(5)}$	0	0	0	0

\tilde{x}	\hat{y}_1	\hat{y}_2
-------------	-------------	-------------

Classifier Chains (and Greedy Inference)

A **classifier chain** is a Bayesian network, with a fully-cascaded structure across output variables. Each $P(Y_j|\text{pa}(Y_j))$ is provided by a **base classifier** (e.g., logistic regression, decision trees, neural network, ...). Motivation: Model **label dependence** with structure.



X	Y_1	Y_2	Y_3	Y_4
$x^{(1)}$	0	1	1	1
$x^{(2)}$	1	0	0	0
$x^{(3)}$	0	1	0	1
$x^{(4)}$	1	0	0	0
$x^{(5)}$	0	0	0	0

\tilde{x}	\hat{y}_1	\hat{y}_2	\hat{y}_3	
			✓	

For example,

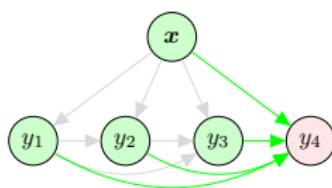
$$\hat{y}_3 = h_3(\mathbf{x}, \hat{y}_1, \hat{y}_2) = \underset{y_3 \in \{0,1\}}{\operatorname{argmax}} P(y_3 | \mathbf{x}, \hat{y}_1, \hat{y}_2)$$

Use training data to fit **base classifier** h_3 .

Inference: $\hat{y}_1, \hat{y}_2, \dots$ are **greedy predictions** from h_1, h_2, \dots

Classifier Chains (and Greedy Inference)

A **classifier chain** is a Bayesian network, with a fully-cascaded structure across output variables. Each $P(Y_j|\text{pa}(Y_j))$ is provided by a **base classifier** (e.g., logistic regression, decision trees, neural network, ...). Motivation: Model **label dependence** with structure.



\mathbf{x}	Y_1	Y_2	Y_3	Y_4
$x^{(1)}$	0	1	1	1
$x^{(2)}$	1	0	0	0
$x^{(3)}$	0	1	0	1
$x^{(4)}$	1	0	0	0
$x^{(5)}$	0	0	0	0

$\tilde{\mathbf{x}}$	\hat{y}_1	\hat{y}_2	\hat{y}_3	\hat{y}_4

For example,

$$\hat{y}_3 = h_3(\mathbf{x}, \hat{y}_1, \hat{y}_2) = \underset{y_3 \in \{0,1\}}{\operatorname{argmax}} P(y_3 | \mathbf{x}, \hat{y}_1, \hat{y}_2)$$

Use training data to fit **base classifier** h_3 .

Inference: $\hat{y}_1, \hat{y}_2, \dots$ are **greedy predictions** from h_1, h_2, \dots

Recap: At **training time**, a multi-label dataset is transformed into L standard binary classification problems:

\mathbf{X}	Y_1	Y_2	Y_3	Y_4
$\mathbf{x}^{(1)}$	0	1	1	0
$\mathbf{x}^{(2)}$	1	0	0	0
$\mathbf{x}^{(3)}$	0	1	0	0
$\mathbf{x}^{(4)}$	1	0	0	1
$\mathbf{x}^{(5)}$	0	0	0	1
$\tilde{\mathbf{x}}$?	?	?	?

At **test time**, each **base classifier** h_j provides $\hat{y}_j \in \{0, 1\}$:

$$\hat{y}_1 = h_1(\tilde{\mathbf{x}}) \quad \triangleright \text{for } j = 1$$

$$\hat{y}_j = h_j(\tilde{\mathbf{x}}, \hat{y}_1, \dots, \hat{y}_{j-1}) \quad \triangleright \text{for } j = 2, \dots, L$$

Multi-output classification:

$$\hat{\mathbf{y}} = [\hat{y}_1, \dots, \hat{y}_L]$$

Recap: At **training time**, a multi-label dataset is transformed into L standard binary classification problems:

$\tilde{\mathbf{x}}$	\hat{y}_1	$\tilde{\mathbf{x}}$	\hat{y}_1	\hat{y}_2	$\tilde{\mathbf{x}}$	\hat{y}_1	\hat{y}_2	\hat{y}_3	$\tilde{\mathbf{x}}$	\hat{y}_1	\hat{y}_2	\hat{y}_3	\hat{y}_4
$x^{(1)}$	0	$x^{(1)}$	0	1	$x^{(1)}$	0	1	1	$x^{(1)}$	0	1	1	0
$x^{(2)}$	1	$x^{(2)}$	1	0	$x^{(2)}$	1	0	0	$x^{(2)}$	1	0	0	0
$x^{(3)}$	0	$x^{(3)}$	0	1	$x^{(3)}$	0	1	0	$x^{(3)}$	0	1	0	0
$x^{(4)}$	1	$x^{(4)}$	1	0	$x^{(4)}$	1	0	0	$x^{(4)}$	1	0	0	1
$x^{(5)}$	0	$x^{(5)}$	0	0	$x^{(5)}$	0	0	0	$x^{(5)}$	0	0	0	1

At **test time**, each **base classifier** h_j provides $\hat{y}_j \in \{0, 1\}$:

$$\hat{y}_1 = h_1(\tilde{\mathbf{x}}) \quad \triangleright \text{for } j = 1$$

$$\hat{y}_j = h_j(\tilde{\mathbf{x}}, \hat{y}_1, \dots, \hat{y}_{j-1}) \quad \triangleright \text{for } j = 2, \dots, L$$

Multi-output classification:

$$\hat{\mathbf{y}} = [\hat{y}_1, \dots, \hat{y}_L]$$

Classifier Chains (Bayes Optimal Inference)

Each base classifier has a probabilistic interpretation:

$$\hat{y}_j = h_j(\mathbf{x}) = \operatorname{argmax}_{\mathbf{y}_j \in \{0,1\}} P(\mathbf{y}_j | \mathbf{x}, y_1, \dots, y_{j-1})$$

(e.g., logistic regression).

From Bayesian network:

$$\begin{aligned}\hat{\mathbf{y}} &= \operatorname{argmax}_{\mathbf{y} \in \{0,1\}^L} P(\mathbf{y} | \mathbf{x}) \\ &= \operatorname{argmax}_{\mathbf{y} \in \{0,1\}^L} P(y_1 | \mathbf{x}) \prod_{j=2}^L P(y_j | \mathbf{x}, y_1, \dots, y_{j-1})\end{aligned}$$

This is not the same as greedy inference.

Classifier Chains (Bayes Optimal Inference)

Each base classifier has a probabilistic interpretation:

$$\hat{y}_j = h_j(\mathbf{x}) = \operatorname{argmax}_{\mathbf{y}_j \in \{0,1\}} P(\mathbf{y}_j | \mathbf{x}, y_1, \dots, y_{j-1})$$

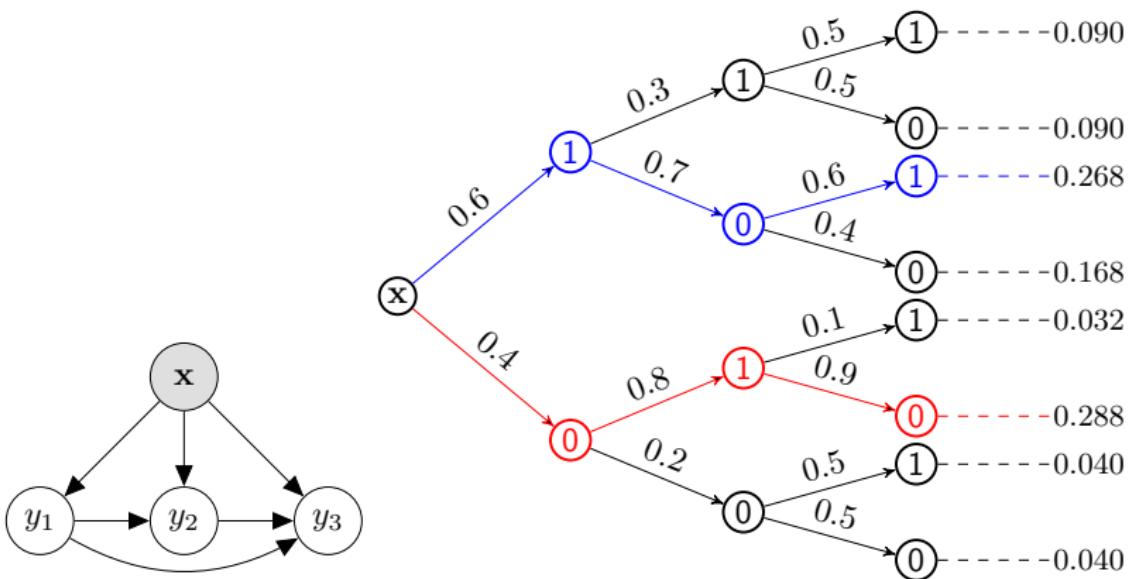
(e.g., logistic regression).

From Bayesian network:

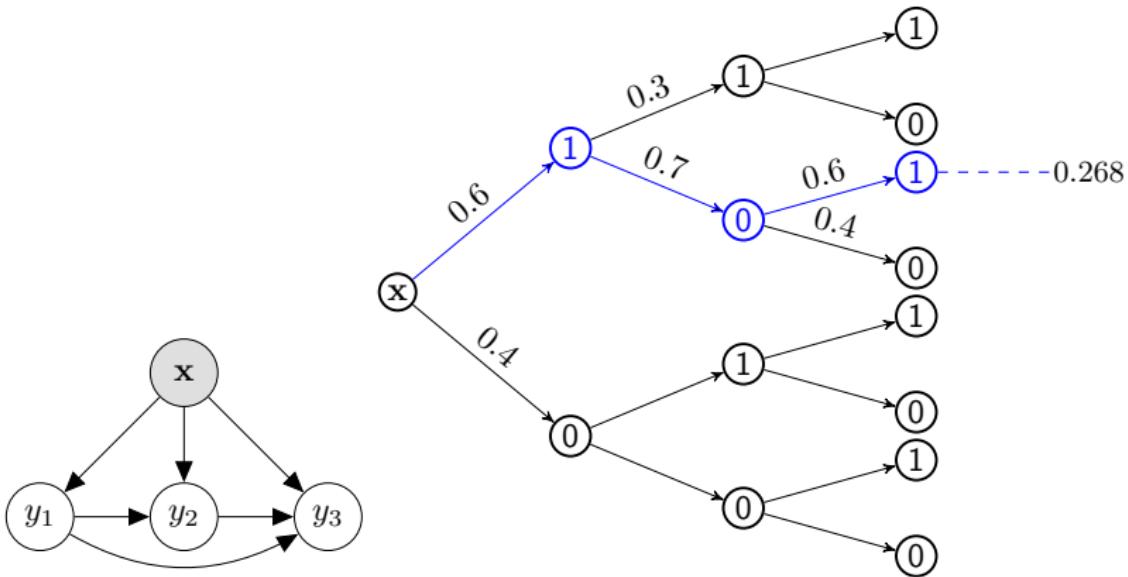
$$\begin{aligned}\hat{\mathbf{y}} &= \operatorname{argmax}_{\mathbf{y} \in \{0,1\}^L} P(\mathbf{y} | \mathbf{x}) \\ &= \operatorname{argmax}_{\mathbf{y} \in \{0,1\}^L} P(y_1 | \mathbf{x}) \prod_{j=2}^L P(y_j | \mathbf{x}, y_1, \dots, y_{j-1})\end{aligned}$$

This is not the same as greedy inference.

First problem: Exponential complexity.

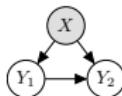


- Exhaustive search $\in \{0, 1\}^L$, i.e., 2^L paths,
Bayes optimal, usually **intractable**
- Single pass down chain (**greedy**), 1 path,
Searches greedily for the joint mode; **not optimal**.



- Exhaustive search $\in \{0, 1\}^L$, i.e., 2^L paths,
Bayes optimal, usually **intractable**
- Single pass down chain (**greedy**), 1 path,
Searches greedily for the joint mode; **not optimal**.

First Question: Is Binary Relevance Enough?



Suppose we want to predict and evaluate in particular

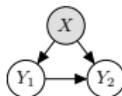
$$\ell(y_2, \hat{y}_2) = \begin{cases} 1 & y_2 \neq \hat{y}_2, \\ 0 & y_2 = \hat{y}_2 \end{cases}$$

\hat{y}_1 does not appear in our query, nor as evidence:

$$P(y_2|\mathbf{x}) = \sum_{y_1 \in \{0,1\}} P(y_1|\mathbf{x})P(y_2|\mathbf{x}, y_1)$$

i.e., we can model $P(y_2|\mathbf{x})$ however we want, directly from \mathbf{x} .

First Question: Is Binary Relevance Enough?



Suppose we want to predict and evaluate in particular

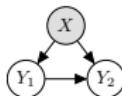
$$\ell(y_2, \hat{y}_2) = \begin{cases} 1 & y_2 \neq \hat{y}_2, \\ 0 & y_2 = \hat{y}_2 \end{cases}$$

\hat{y}_1 does not appear in our query, nor as evidence:

$$P(y_2|\mathbf{x}) = \sum_{y_1 \in \{0,1\}} P(y_1|\mathbf{x})P(y_2|\mathbf{x}, y_1)$$

i.e., we can model $P(y_2|\mathbf{x})$ however we want, directly from \mathbf{x} .

First Question: Is Binary Relevance Enough?



Suppose we want to predict and evaluate in particular

$$\ell(y_2, \hat{y}_2) = \begin{cases} 1 & y_2 \neq \hat{y}_2, \\ 0 & y_2 = \hat{y}_2 \end{cases}$$

\hat{y}_1 does not appear in our query, nor as evidence:

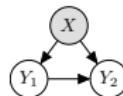
$$P(y_2|\mathbf{x}) = \sum_{y_1 \in \{0,1\}} P(y_1|\mathbf{x})P(y_2|\mathbf{x}, y_1)$$

i.e., we can model $P(y_2|\mathbf{x})$ however we want, directly from \mathbf{x} .

Now suppose we are minimizing **subset 0/1 loss**,

$$\ell(\mathbf{y}, \hat{\mathbf{y}}) = \begin{cases} 1 & \mathbf{y} \neq \hat{\mathbf{y}}, \\ 0 & \mathbf{y} = \hat{\mathbf{y}} \quad (\text{exactly}) \end{cases}$$

First Question: Is Binary Relevance Enough?



Suppose we want to predict and evaluate in particular

$$\ell(y_2, \hat{y}_2) = \begin{cases} 1 & y_2 \neq \hat{y}_2, \\ 0 & y_2 = \hat{y}_2 \end{cases}$$

\hat{y}_1 does not appear in our query, nor as evidence:

$$P(y_2|\mathbf{x}) = \sum_{y_1 \in \{0,1\}} P(y_1|\mathbf{x})P(y_2|\mathbf{x}, y_1)$$

i.e., we can model $P(y_2|\mathbf{x})$ however we want, directly from \mathbf{x} .

Now suppose we are minimizing **subset 0/1 loss**,

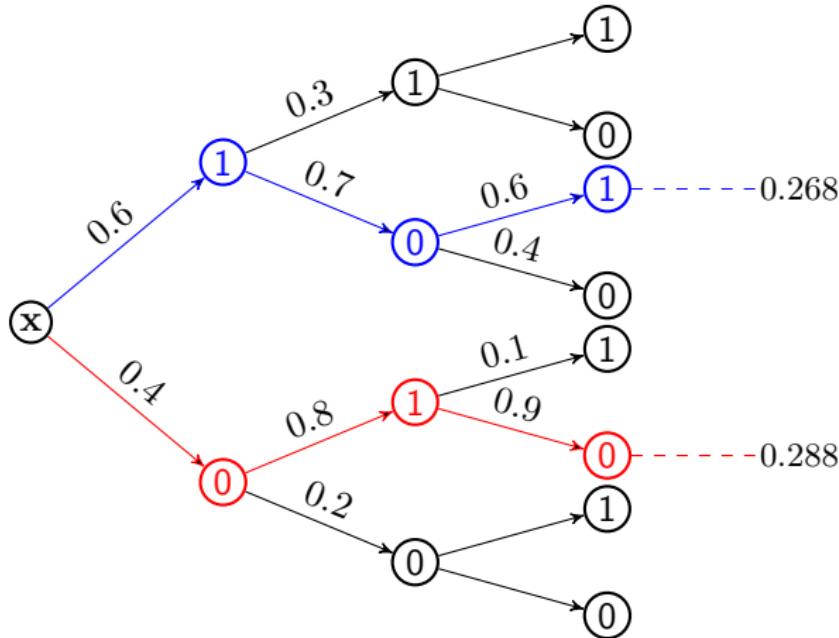
$$\ell(\mathbf{y}, \hat{\mathbf{y}}) = \begin{cases} 1 & \mathbf{y} \neq \hat{\mathbf{y}}, \\ 0 & \mathbf{y} = \hat{\mathbf{y}} \quad (\text{exactly}) \end{cases}$$

Now suppose mean squared error $MSE(\mathbf{y}, \hat{\mathbf{y}})$.

Inference as Search

- 1 Multi-Label Classification
- 2 Classifier Chains
- 3 Inference as Search
- 4 Benefits of Structure
- 5 Structure from Data / Learning Structure as a Search
- 6 Summary

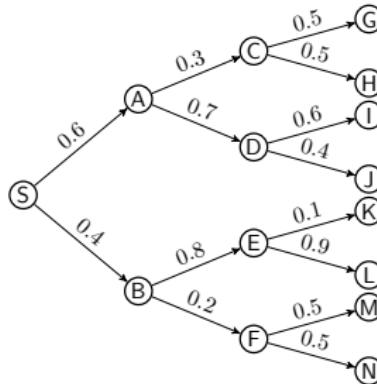
Inference as Tree-Search



- Node $n := y_1, \dots, y_j \in \{0, 1\}^j$; $j \in \{1, \dots, L\}$ and start state \emptyset
- Goal state: Some $\mathbf{y} = y_1, \dots, y_L$
- Edge cost: $P(y_j | \mathbf{x}, y_1, \dots, y_{j-1})$
- Path cost: $g^{-1}(n) = P(y_1, \dots, y_j | \mathbf{x})$. • N.B. \neq num. of steps!

Uniform Cost Search¹ (UCS)

Expanding node n will incur a **path-cost** $g(n)$.



- A **priority queue** of nodes and respective **path cost** $\{(n, g(n))\}$
- Expand lower-cost nodes first (n giving smallest $g(n)$)
- Goal test is applied *when selected for expansion* rather than first encountered (in case a better path there can be found)
- **Optimal** in general (for costs ≥ 0), i.e., the solution is the best

¹aka **Dijkstra's Algorithm**, esp. when expanding *all* nodes

On weighted graph G , from node n_0 to (goal node) n_{goal} :

UNIFORMCOSTSEARCH(G, n_0, n_{goal}):

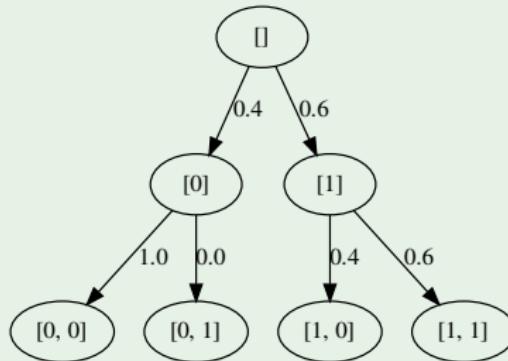
- $\mathcal{Q} = \{(n_0, g(n_0))\}$ ▷ The priority queue
- While $|\mathcal{Q}| > 0$:
 - $n \leftarrow \mathcal{Q}.\text{pop}()$ ▷ Pop² n with smallest $g(n)$
 - if $n = n_{\text{goal}}$, stop!
 - for $n' \in \text{Neighbours}(n)$:
 - $g(n') = g(n) + \text{cost}(n, n')$
 - $\mathcal{Q} \leftarrow \mathcal{Q} \cup (n', g(n'))$

(N.B. Usually also want to keep track of the path).

²where

$$n^* \leftarrow \mathcal{Q}.\text{pop}() \Leftrightarrow \begin{cases} n^* = \underset{n \in \mathcal{Q}}{\text{argmin}} g(n), \\ \mathcal{Q} \leftarrow \mathcal{Q} \setminus n^* \end{cases}$$

Uniform Cost Search for Probability-Tree Search



Using queue \mathcal{Q} , of tuples $(n, g(n))$,

$$\mathcal{Q} = \{([], 0)\}$$

$$\mathcal{Q} = \{([1], 0.6), ([0], 0.4)\} \quad \triangleright \text{pop } []$$

$$\mathcal{Q} = \{([0], 0.4), ([1, 1], 0.36), ([1, 0], 0.24)\} \quad \triangleright \text{pop } [1]$$

$$\mathcal{Q} = \{([0, 0], 0.4), ([1, 1], 0.36), ([1, 0], 0.24), ([0, 1], 0.0)\} \quad \triangleright \text{pop } [0]$$

Return: $y = [0, 0]$

Beam Search

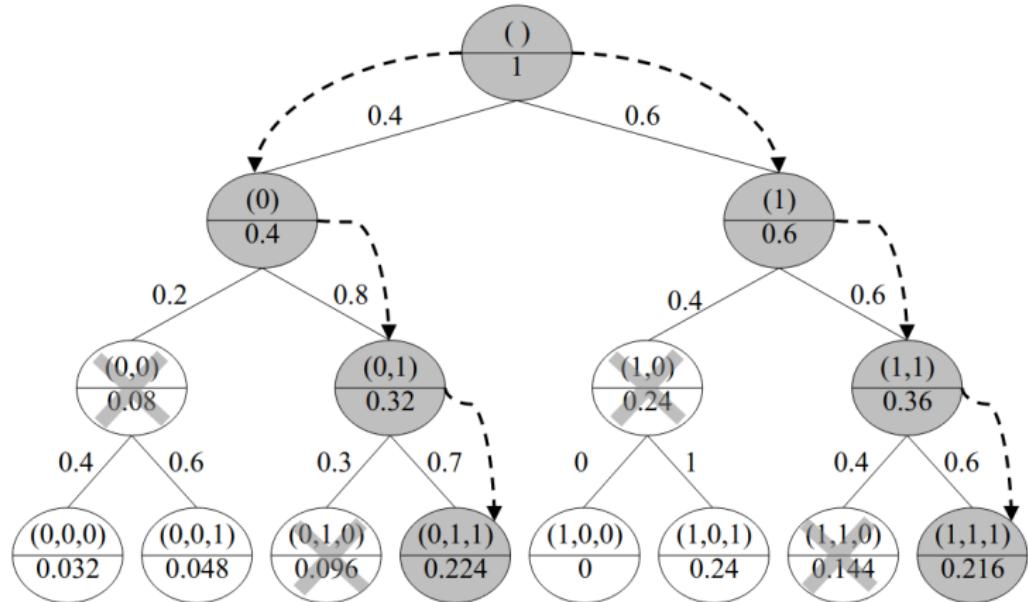
A variant of UCS with reduced memory requirements; only considers the top β nodes (the size of the **beam**) wrt path cost at each level.

- A variant of greedy search.
- Popular in neural networks, especially in language models

Where β is a hyperparameter;

- If $\beta = \infty$, then it is standard UCS

Beam Search, $\beta = 2$



From: Mena et al., *An overview of Inference Methods in Probabilistic Classifier Chains for Multi-label classification*

ϵ -Approximate Search

Proceed like UCS (node $[y_1, \dots, y_j]$ is added to the queue), *but only if*, at depth j ,

$$\underbrace{P(y_1, \dots, y_j | x)}_{g(n)} > \epsilon$$

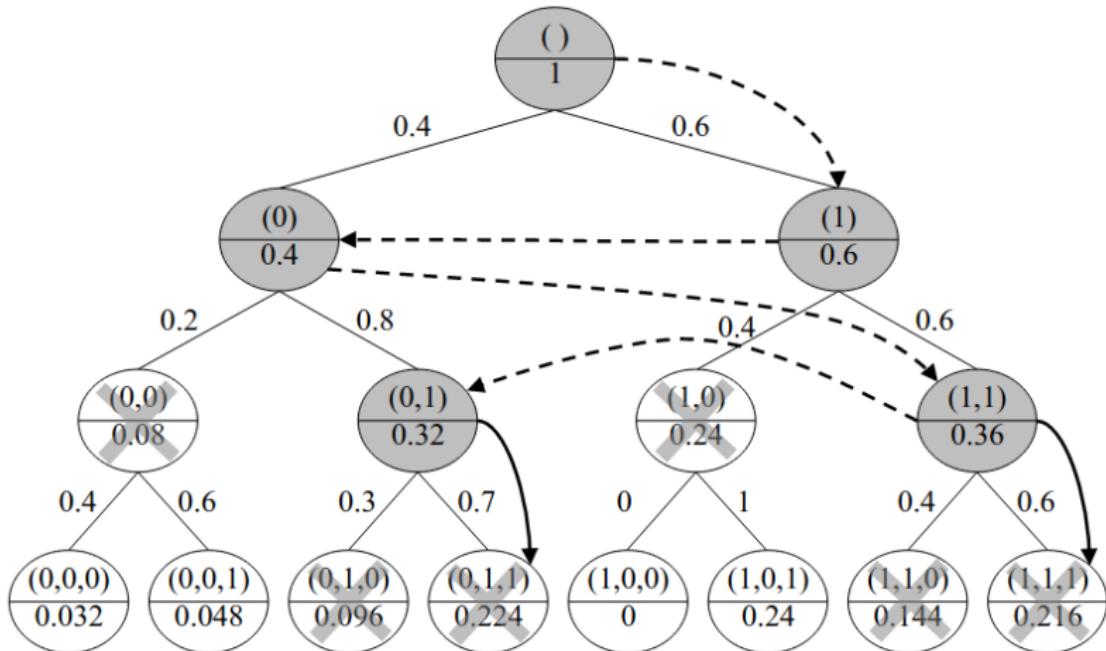
Continue (pop, goal-check, add children), until

- a leaf is found (done!); or
- queue is empty: apply **greedy search** to goal
(from all childless (non-expanded) nodes).

Where ϵ is a hyper-parameter;

- if $\epsilon = 0.0$, then unit costs are ignored; defaults to **UCS**
- if $\epsilon = 0.5$, defaults to **greedy search**.

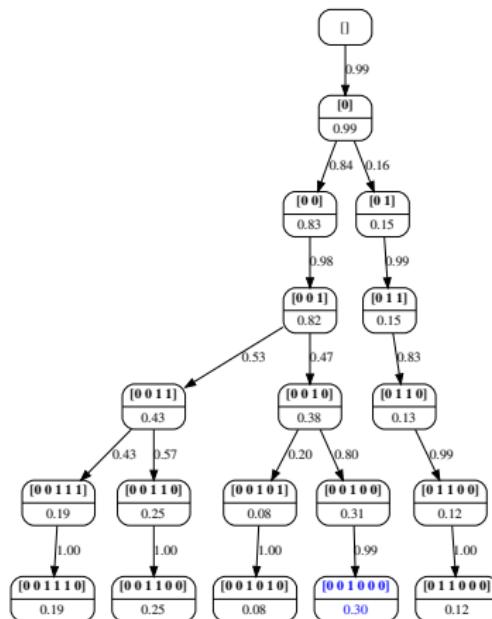
ϵ -Approximate Search, $\epsilon = 0.25$.



From: Mena et al., An overview of Inference Methods in Probabilistic Classifier Chains for Multi-label classification

Monte Carlo Search

- ① Take M samples/**paths** from the search tree from root to leaf i.e., open branch $y_j^{(m)}$ with prob. $P(y_j^{(m)} | \mathbf{x}, y_1^{(m)}, \dots, y_{j-1}^{(m)})$
- ② Return the m^* -th path; where $m^* = \operatorname{argmax}_m P(\mathbf{y}^{(m)} | \mathbf{x})$.



On the Emotions/Music dataset: $M = 20$, only explored branches shown

Recall: We were optimizing **0/1 loss**.

- Chance of finding Bayes-optimal solution depends on $\beta/\epsilon/M$
- Other search methods are applicable³, e.g., A* search,

³Survey:

Recall: We were optimizing **0/1 loss**.

- Chance of finding Bayes-optimal solution depends on $\beta/\epsilon/M$
- Other search methods are applicable³, e.g., A* search,

The End ?

But – Empirical observation: Classifier chains also outperforms baseline methods under **Hamming loss** (where **labels evaluated independently**); a decomposable metric – risk minimization tells us that classifier chains provides no advantage).

³Survey:

Benefits of Structure

- 1 Multi-Label Classification
- 2 Classifier Chains
- 3 Inference as Search
- 4 Benefits of Structure
- 5 Structure from Data / Learning Structure as a Search
- 6 Summary

Why Model Labels Together

Other reasons to model labels together (other than ‘because of **label dependence**/to minimise 0/1-loss’):

- Connectivity \Rightarrow greater **efficiency** (sometimes)
- Connectivity \Rightarrow better **interpretation**
- Connectivity \Rightarrow **regularization**
- Connectivity \Rightarrow **predictive power** (deep learning!)
- Connectivity \Rightarrow suited to more loss metrics

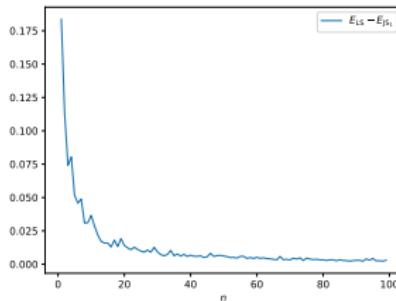
Reasons **not** to model targets together: greater **efficiency** (sometimes), **simplicity**, **parallelizability**, ...

Reason #1: Regularization

Joint-target regularization (modeling labels together) is beneficial even if targets are intrinsically independent.

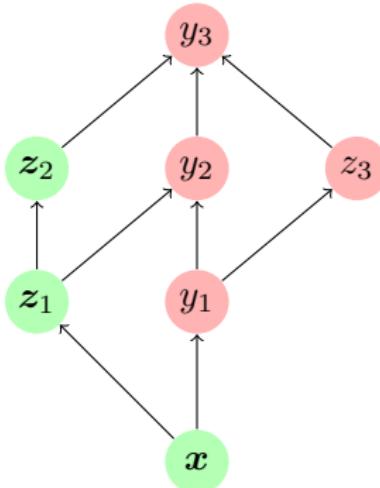
James-Stein Estimator (m targets, n samples)

$$\begin{aligned}\hat{\mathbf{y}}_{JS_n} &= \frac{1 - (m - 2)\frac{\hat{\sigma}^2}{n}}{\|\bar{\mathbf{y}}\|^2} \cdot \bar{\mathbf{y}} \\ &= \lambda(\{\mathbf{y}^{(1)}, \dots, \mathbf{y}^{(n)}\}) \cdot \hat{\mathbf{y}}_{MLE}\end{aligned}$$



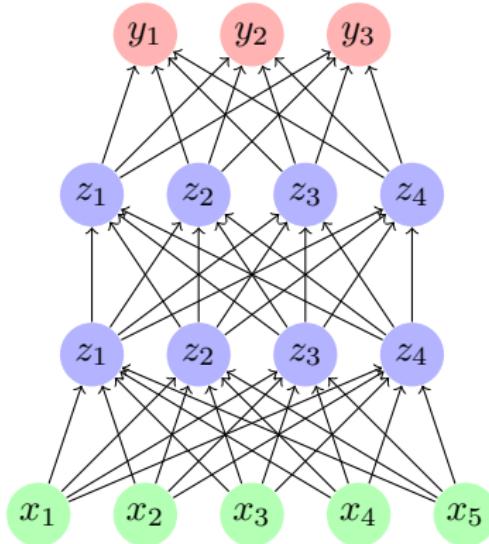
But the advantages quickly fade as $n \gg 1$.

Reason #2: Predictive Power/Architecture

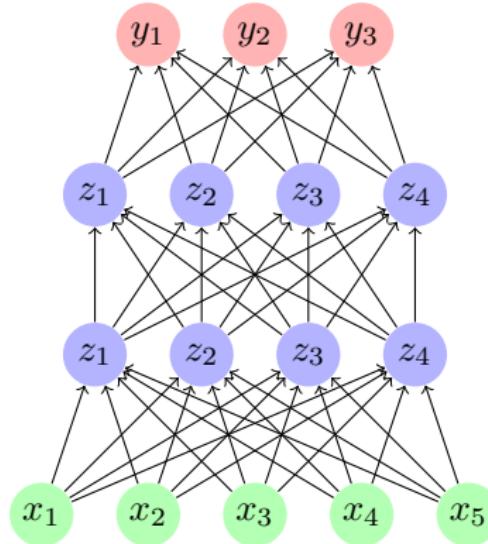


with delay nodes $z = f(x) = x$

- Forward propagation = greedy inference
- It's **deep in the label space!**
- labels = feature space/hidden nodes 'for free'
- More parameters \Rightarrow more powerful model



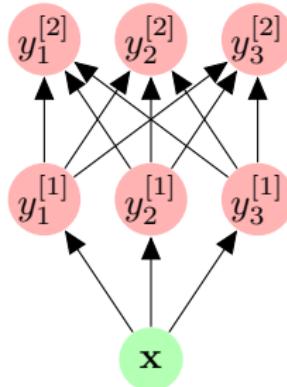
(This also works, but we have to learn $z^{[l]}$ from scratch.)



(This also works, but we have to learn $z^{[l]}$ from scratch.)
There are ‘deep’ chaining-type architectures; consider, e.g., ResNet, skip layers, etc.

Correcting Errors/Bias with Stacking

In multi-label stacking, the predictions are stacked as features.



$$y_j^{[l]} = j\text{-th label, } l\text{-th layer}$$

Like classifier chains at inference time: $\hat{y}_j^{[1]}$ is used to predict $\hat{y}_k^{[2]}$,
but important **difference in training**: predictions $\hat{y}_j^{(i)}$ are used!

- Using dataset labels: model conditional dependencies, better under 0/1 loss
- Using predicted labels: extend feature space, improve under Hamming loss

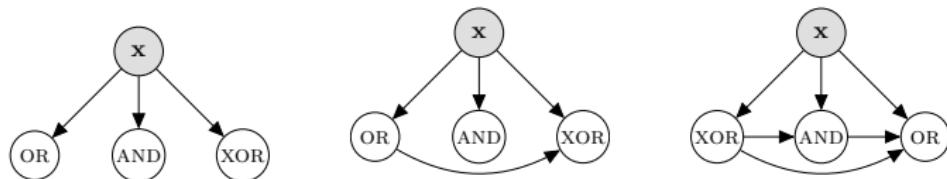
Comparison

	X_1	X_2	X_3	Y_2	X_1	X_2	X_3	Y_2
Basis expansion	x	ϕ_1	ϕ_2	y_2	\tilde{x}	ϕ_1	ϕ_2	\hat{y}_2
Classifier chain	x	y_1		y_2	\tilde{x}	\hat{y}_1		\hat{y}_2
Stacking	x	\hat{y}_1	\hat{y}_2	y_2	\tilde{x}	$\hat{y}_1^{[1]}$	$\hat{y}_2^{[1]}$	$\hat{y}_2^{[2]}$
Neural network	x			y_2	\tilde{x}	\hat{z}_1	\hat{z}_2	\hat{y}_2

Training (left; $x \equiv x^{(i)}$, $y_j \equiv y_j^{(i)}$) vs Testing (right; $\hat{y}_2|\tilde{x}$)

Which Structure?

OK, yes, structure works – but which structure?



$x \in \{0, 1\}^2$, and labels are logical operations

Consider:

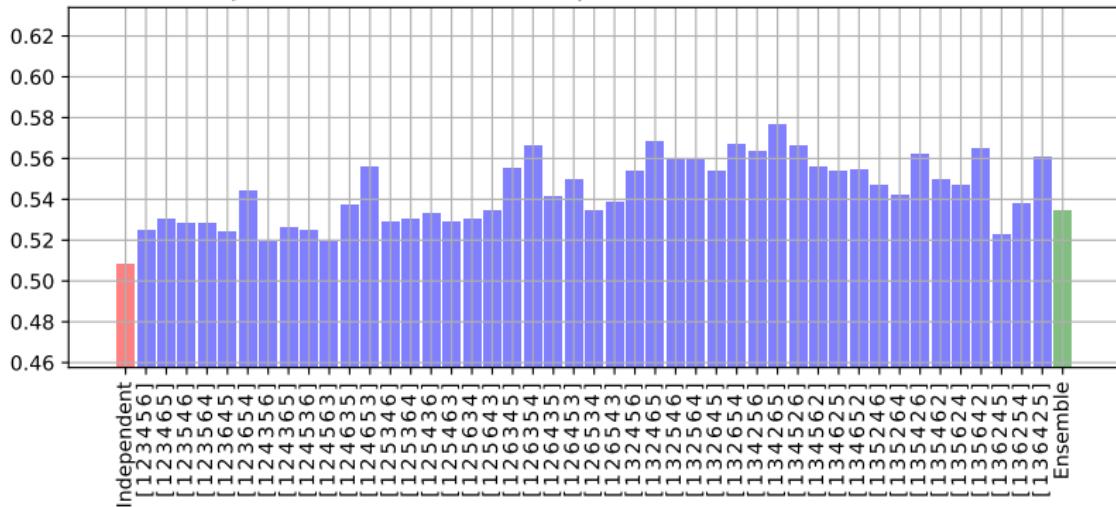
- Choice of base classifier
- Choice of inference
- What is your loss metric?
- Will we find ground truth? Do we need to?

Structure from Data / Learning Structure as a Search

- 1 Multi-Label Classification
- 2 Classifier Chains
- 3 Inference as Search
- 4 Benefits of Structure
- 5 Structure from Data / Learning Structure as a Search
- 6 Summary

Motivation: Structure Can Make a Difference

Jaccard score from 45 chain permutations; 'emotions' data.

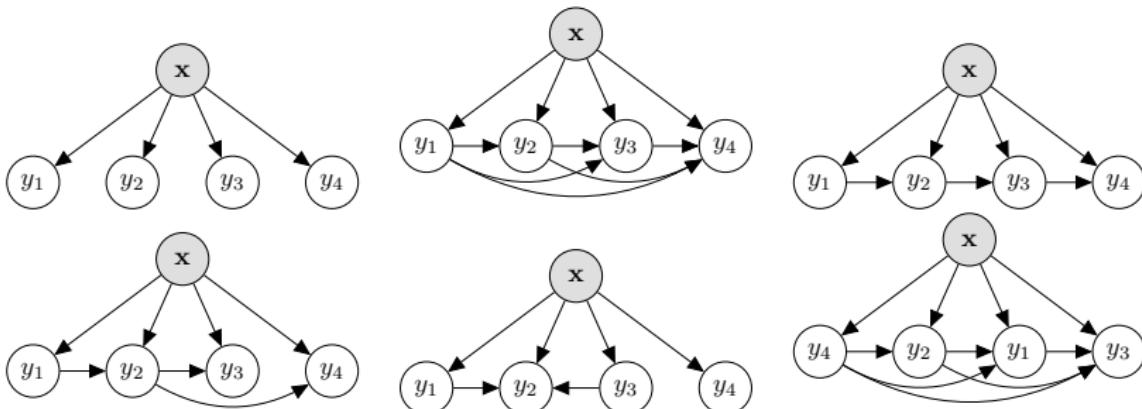


Highlighting some example cases – What is going on?

m	Chain Order	Jaccard	A1	A2	A3	A4	A5	A6
5	[1 2 3 6 4 5]	0.52	0.81	0.70	0.73	0.93	0.82	0.74
6	[1 2 3 6 5 4]	0.54	0.81	0.70	0.73	0.93	0.86	0.74
17	[1 2 5 4 6 3]	0.53	0.81	0.70	0.73	0.94	0.81	0.73
18	[1 2 5 6 3 4]	0.53	0.81	0.70	0.71	0.91	0.81	0.77
20	[1 2 6 3 5 4]	0.57	0.81	0.70	0.75	0.95	0.84	0.77
42	[1 3 5 6 4 2]	0.57	0.81	0.72	0.74	0.94	0.83	0.76

Jaccard score and per-label accuracy (A_j) for all labels $j = 1, \dots, 6$

Options for Structure



- ➊ Random structure (often in an ensemble).
- ➋ Use an existing hierarchy (expert knowledge)
- ➌ Use a full/complete structure
- ➍ Impose a fixed structure allowing tractable inference
- ➎ Learning a structure via procedure/[search](#) based on a score of
 - ➏ Heuristics based on label dependence,
 - ➏ Individual accuracy
 - ➏ Pairwise accuracy
 - ➏ Structure accuracy
 - ➏ ...

Learning Structure from Data

Assume training **data**

$$\mathcal{D} \sim P_{G^*}$$

(is generated from ground-truth **graph** $G^* \in \mathcal{G}$).

Unfortunately:

- Intractable number of graphs $|\mathcal{G}|$ (super-exponential)
- And we cannot identify G^* anyway
 - Not enough data; instability (different G for same \mathcal{D})
 - Different G for same P (expression of joint dist.)

Fortunately:

- We can recover some equivalent G / equivalence class of G^*
- We often don't need G^* (just some G for good accuracy)
- Many approaches are available to help us **search** for it

Search for Structure

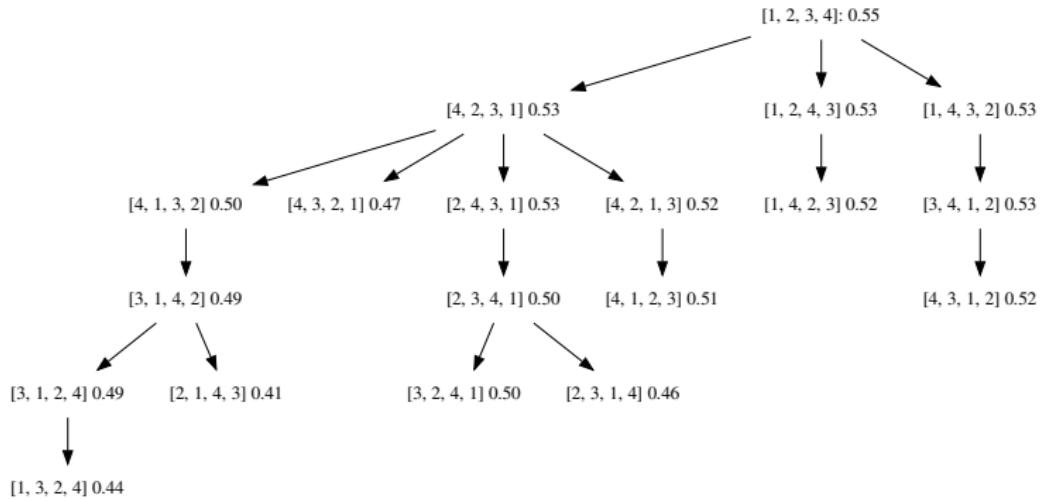
- Search space: \mathcal{G} (all possible graphs/structures)
- Nodes n (let n identify a graph $G_n \in \mathcal{G}$)
- Node cost $g(n) \approx$ expected loss incurred by our model h_{G_n}

Our search problem:

$$\hat{G} = \operatorname{argmin}_{G_n \in \mathcal{G}} \underbrace{\mathbb{E}[\ell(h_{G_n}(\mathbf{X}), \mathbf{Y})]}_{\approx g(n)}$$

- We can apply the same strategies used for inference (UCS, Beam-Search, Monte Carlo search, ...) and others.
- Expensive (must approximate \mathbb{E} , implies building models)
- But there are tricks, e.g., search only for order, freeze the graph from root to leaf during the search.

Remark/reminder: consideration of our cost function $g(n)$, ℓ is important!



Dynamic Structures

So far:

- ① Find $\hat{G} \in \mathcal{G}$, train $h_{\hat{G}}$
- ② At inference time,

$$\hat{\mathbf{y}} = h_{\hat{G}}(\tilde{\mathbf{x}}) = \operatorname{argmax}_{\mathbf{y} \in \mathcal{Y}} P_{\hat{G}}(\mathbf{y} | \tilde{\mathbf{x}})$$

but is \hat{G} best for test $\tilde{\mathbf{x}}$ in particular?

- Structure may be conditionally dependent on test instance!
- We can include structure in inference;

$$\hat{\mathbf{y}} = h(\tilde{\mathbf{x}}) = \operatorname{argmax}_{\mathbf{y} \in \mathcal{Y}, G \in \mathcal{G}} P_G(\mathbf{y} | \mathbf{x})$$

approximations are necessary; many trade-offs involved!

Ideal setting: an **agent** looks at the data, provides the appropriate structure and parameters (architecture).

Neural Architecture Search

- Structure search is not specific to Bayesian networks
- Neural networks face the same problem: how to know what structure (architecture) and [hyper]-parameters to use, for a given problem/data set
- Bayesian networks can be neural networks (e.g., classifier chains); (Generally: Bayesian networks marginalize out z , Neural networks propagate through z)

Some solutions/related concepts to be aware of:

- Use a **recurrent structure** (e.g., RNNs, LSTMs; output labels as a sequence)
- Transformers and **attention** mechanisms
- Hyperparameter optimization (**structure as a hyperparameter**)
- **Meta learning**: learning to learn (including structure)
- Automated ML (**AutoML**): structure choice is automated
- **Neural Architecture Search**: a **policy** to decide how build and traverse structure; **reinforcement learning**.

Summary

- 1 Multi-Label Classification
- 2 Classifier Chains
- 3 Inference as Search
- 4 Benefits of Structure
- 5 Structure from Data / Learning Structure as a Search
- 6 Summary

Summary

Machine learning/AI applications are increasingly multi- and structured-output. Models often require **structure**.

How to address this issue? Take a principled approach! –

- Do you **need** structure? (consider your loss metric)
- What **purpose** will it serve? (architecture for predictive power? space efficiency? regularization? interpretation?)
- Is one structure from training sufficient? Or do you need several structures, dynamically?
- How will you traverse your structure (what type of **inference**)?

Probabilistic vs [deep] neural-network architectures:

- Much in common, can be equivalent/part of the same model
- Different perceptions of “depth”: layers vs iterations

Everything a search: Learning (find $\hat{G}, \hat{\theta}$) and Inference (find $\hat{\mathbf{y}}$):

$$\hat{\mathbf{y}} = \operatorname{argmax}_{\mathbf{y} \in \{0,1\}^L, G_\theta \in \mathcal{S}} P_{G_\theta}(\mathbf{y} | \mathbf{x})$$

Structure benefits agents and agents benefit structure [search]!

Multi-Target and Structured-Output Learning

INF581 Advanced Machine Learning and Autonomous Agents

Jesse Read

